

# Εργασία στα Γραμμικά Μοντέλα

Κωνσταντίνος Παπανίκος (1112201600167)

Κωνσταντίνα Τύραλη (1112201600230)

Δημήτρης Φούντας (1112201600236)

23 Ιουνίου 2020



ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ  
**Εθνικόν και Καποδιστριακόν**  
**Πανεπιστήμιον Αθηνών**  
———ΙΔΡΥΘΕΝ ΤΟ 1837———

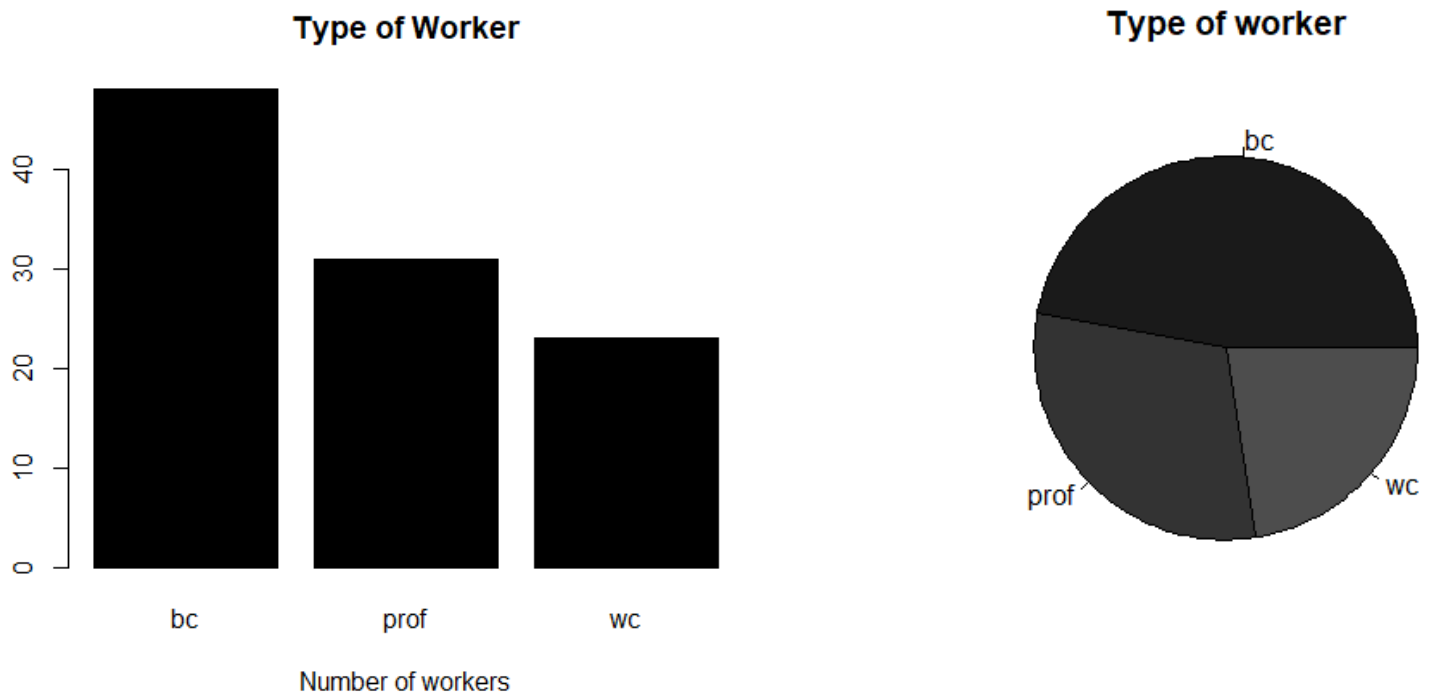
## ΠΕΡΙΛΗΨΗ:

Αρχικά παρατηρούμε ότι το δείγμα μας όσον αφορά τον τύπο εργαζομένων δεν είναι ιδιαίτερα ομοιόμορφο, έχουμε 48 χειρωνακτές, 31 διοικητικούς και μόλις 23 επαγγελματίες γραφείου. Εξίσου εύκολα προσέχουμε πως το μέσο *prestige* του δείγματος μας είναι 46 στην κλίμακα που θέσπισαν οι *Pineo – Porter* με την ψυχοκοινωνική μέθοδο που ανέπτυξαν κατά τη δεκαετία του 1960, η μέση εκπαίδευση είναι 10.54 χρόνια και το μέσο ποσοστό γυναικών ανά επάγγελμα είναι 13,6%. Βλέπουμε επίσης ότι υπάρχουν επαγγέλματα στο δείγμα μας τα οποία απαιτούν βασική μόρφωση 6,38 χρόνων ενώ υπάρχουν και επαγγέλματα που απαιτούν κοντά στα 16 χρόνια εκπαίδευσης. Παρομοίως, παρατηρούμε ότι παρ' όλο που το ποσοστό γυναικών ανά επάγγελμα είναι κατά μέσο όρο 13%, γίνεται εμφανές στα στοιχεία μας ότι υπάρχει τουλάχιστον ένα 'γυναικοκρατούμενο' επάγγελμα, έχοντας ποσοστά γυναικών εργαζομένων στο 97,5% (γραμματείς), ενώ υπάρχουν περισσότερα από ένα επαγγέλματα στα οποία εργάζονται αποκλειστικά άντρες (πυροσβέστες, μηχανικοί τρένων κλπ). Παρατηρούμε επίσης ότι η εργασία με το ελάχιστο εισόδημα είναι οι *babysitters* με αμοιβή 611\$ ενώ το μεγαλύτερο εισόδημα του δείγματος μας βρίσκεται στο επάγγελμα του γενικού διευθυντή και ανέρχεται στα 25879\$. Τέλος, βλέπουμε ότι το μικρότερο *prestige* (14.8) βρέθηκε στο επάγγελμα των *newsboys* ενώ το μεγαλύτερο (87.2) βρέθηκε στο επάγγελμα *physicians*. Στην συνέχεια παρατηρούμε ότι συνήθως όσο μεγαλύτερη είναι η εκπαίδευση ενός ατόμου (σε έτη), τόσο μεγαλύτερο και το *prestige* της εργασίας του. Όμως όταν συγκρίναμε την κοινωνικοοικονομική θέση με το εισόδημα παρατηρήσαμε κάτι εξαιρετικά ενδιαφέρον, όταν το *prestige* είναι μικρότερο του 80, όσο μεγαλώνει αυτό μεγαλώνει ανάλογα και το εισόδημα, αλλά όταν το *prestige* ξεπερνάει το 80 το εισόδημα του επαγγελματία αυξάνεται με πολύ μεγαλύτερους ρυθμούς από πριν, με περαιτέρω αύξηση του *prestige*. Κάταλήξαμε επίσης στο συμπέρασμα ότι η κοινωνικοοικονομική θέση που αντιστοιχεί σε ένα επάγγελμα εξαρτάται άμεσα από το εισόδημα και την εκπαίδευση, μάλιστα αν υπάρχει γνώση και των δύο, μπορεί να γίνει μια αρκετά καλή πρόβλεψη για το πόσο *prestige* έχει το συγκεκριμένο επάγγελμα. Για παράδειγμα, ένας εργαζόμενος με 10 χρόνια εκπαίδευσης και 4000\$ εισόδημα προβλέπεται να έχει περίπου 42.87 μονάδες στην κλίμακα του *prestige*. Στην συνέχεια εξετάζουμε το *prestige* σε σχέση με το εισόδημα και την εκπαίδευση ταυτόχρονα αντί για να τα εξετάσουμε ξεχωριστά και βλέπουμε ότι το μοντέλο αυτό μας δίνει πιο αξιόπιστα αποτελέσματα απ' ότι όταν τα διερευνήσαμε ξεχωριστά. Τέλος, χρειάστηκε έλεγχος με ανάλυση διασποράς για το εισόδημα σε σχέση με το αν το ποσοστό γυναικών εν ενεργεία επαγγελματιών είναι πάνω ή κάτω από 50% και το είδος του επαγγέλματος. Ουσιαστικά αυτό έγινε για να εξετάσουμε αν το είδος του επαγγέλματος και το ποσοστό των γυναικών επηρεάζουν τον αναμενόμενο μισθό. Έτσι δημιουργούνται 6 διαφορετικές περιπτώσεις:

- α)πάνω από 50% γυναίκες και χειρωνακτές → αναμενόμενος μισθός: 3.200\$
- β)πάνω από 50% γυναίκες και επαγγελματίες, διοικητικοί και τεχνικοί → αναμενόμενος μισθός: 6.000\$
- γ)πάνω από 50% γυναίκες και επαγγελματίες γραφείου → αναμενόμενος μισθός: 4.000\$
- δ)κάτω από 50% γυναίκες και χειρωνακτές → αναμενόμενος μισθός 6.000\$
- ε)κάτω από 50% γυναίκες και επαγγελματίες, διοικητικοί και τεχνικοί → αναμενόμενος μισθός: 11.800\$
- στ)κάτω από 50% γυναίκες και επαγγελματίες γραφείου → αναμενόμενος μισθός 6.500\$

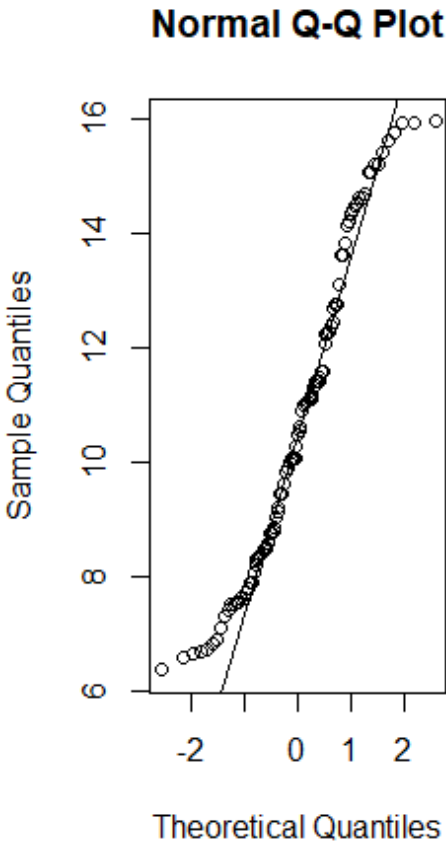
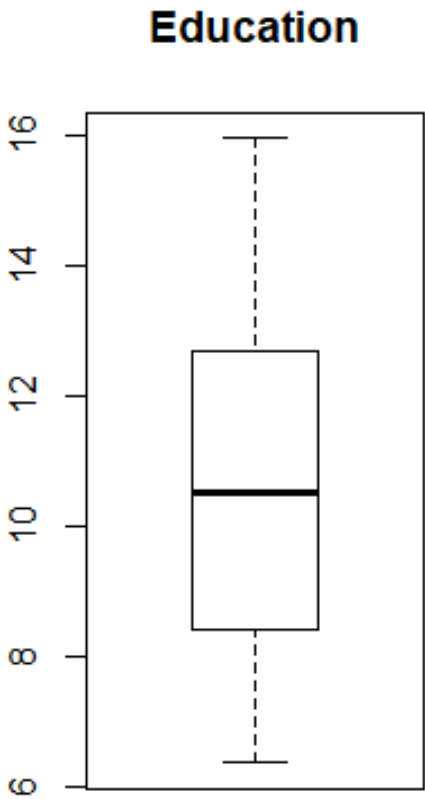
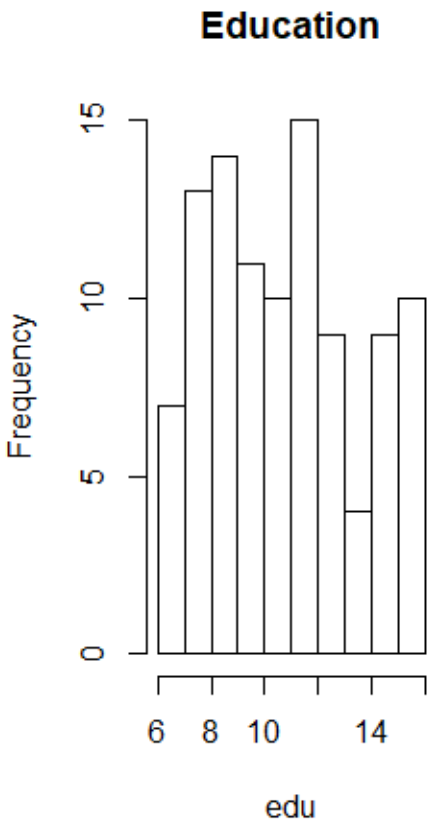
1ο Ερώτημα:

Για αρχή, χωρίζουμε τις μεταβλητές σε δύο κατηγορίες: *qualitative* και *quantitative*, καθώς η κάθε κατηγορία απαιτεί διαφορετική διαχείριση. Η μεταβλητή *type* δεν μπορεί να ποσοτικοποιηθεί, καθώς πρόκειται για τα τρία επίπεδα επαγγελματών, άρα ανήκει στην πρώτη περίπτωση. Έτσι, κάνουμε περιγραφική στατιστική και βρίσκουμε τον αριθμό των εργαζομένων από κάθε επίπεδο: 48 άτομα χειρώναχτες (*bc*), 31 επαγγελματίες, διοικητικοί και τεχνικοί (*prof*) και 23 επαγγελματίες γραφείου (*wc*). Έχουμε τα γραφήματα *barplot* και *pie*:

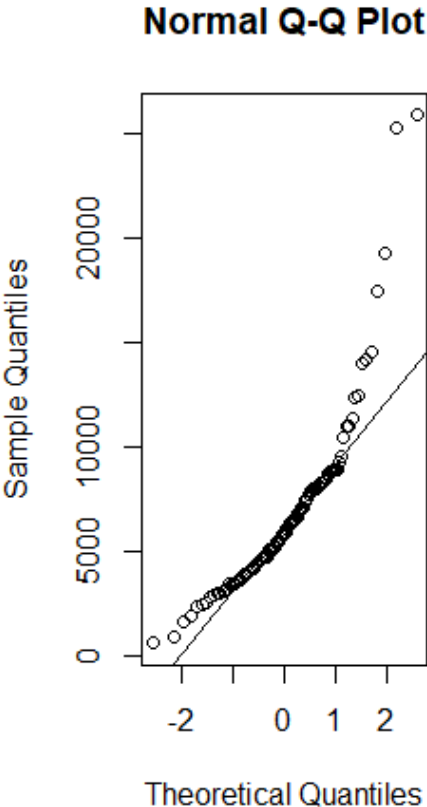
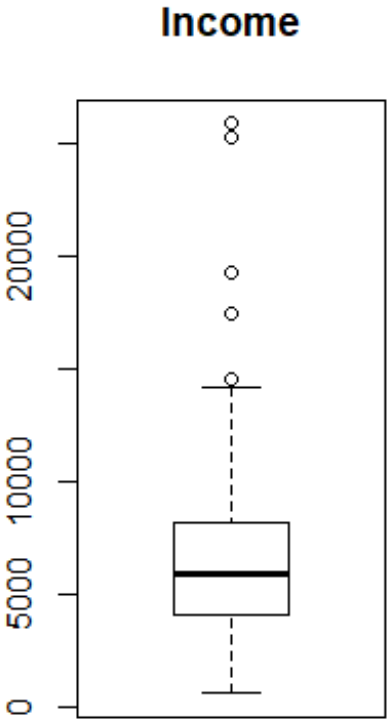
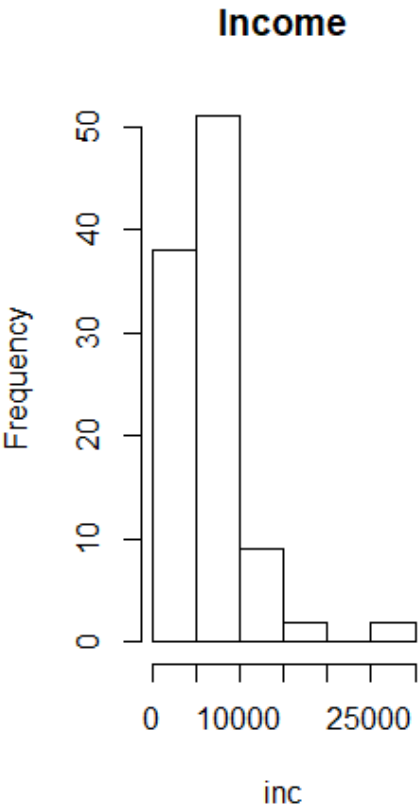


Τώρα, θέλουμε να κάνουμε περιγραφική στατιστική στα υπόλοιπα δεδομένα. Πρώτα θα ελέγξουμε ποιες από τις μεταβλητές μπορεί να ακολουθούν κανονική κατανομή μέσω ιστογράμματος *boxplot* και *q - qplot*. Έχουμε λοιπόν:

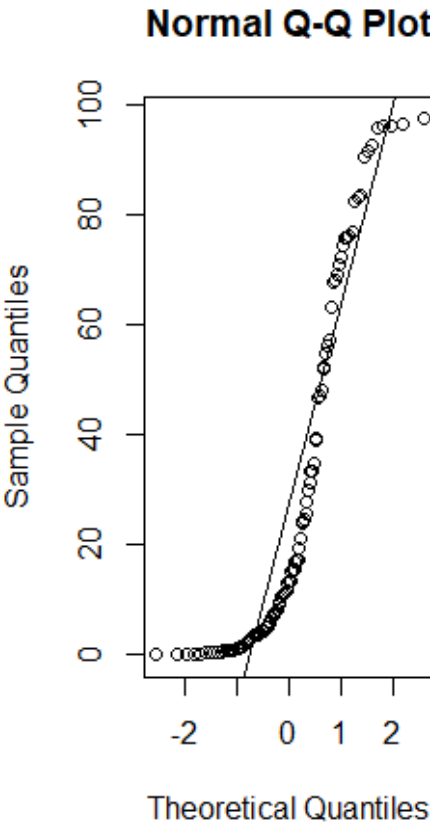
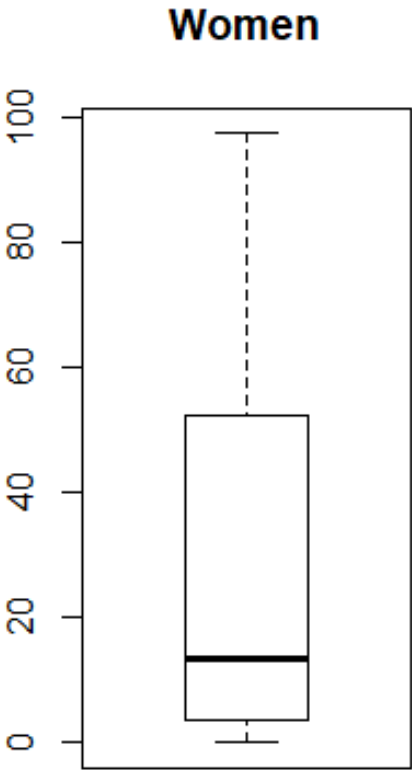
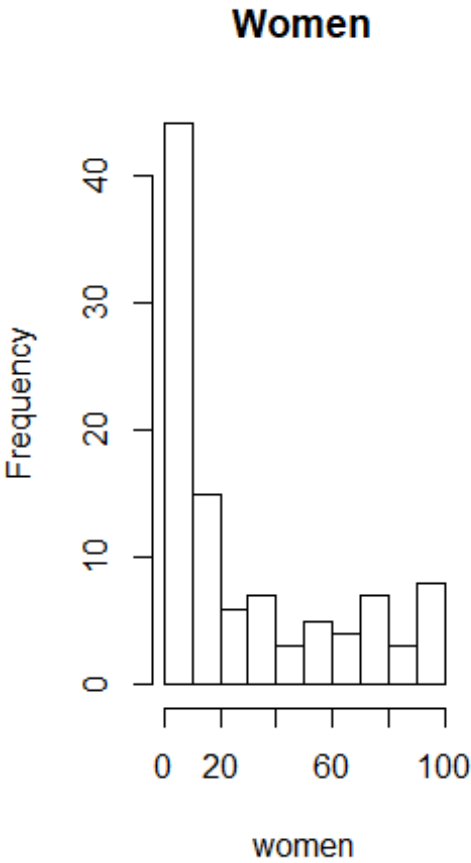
Για την εκπαίδευση:



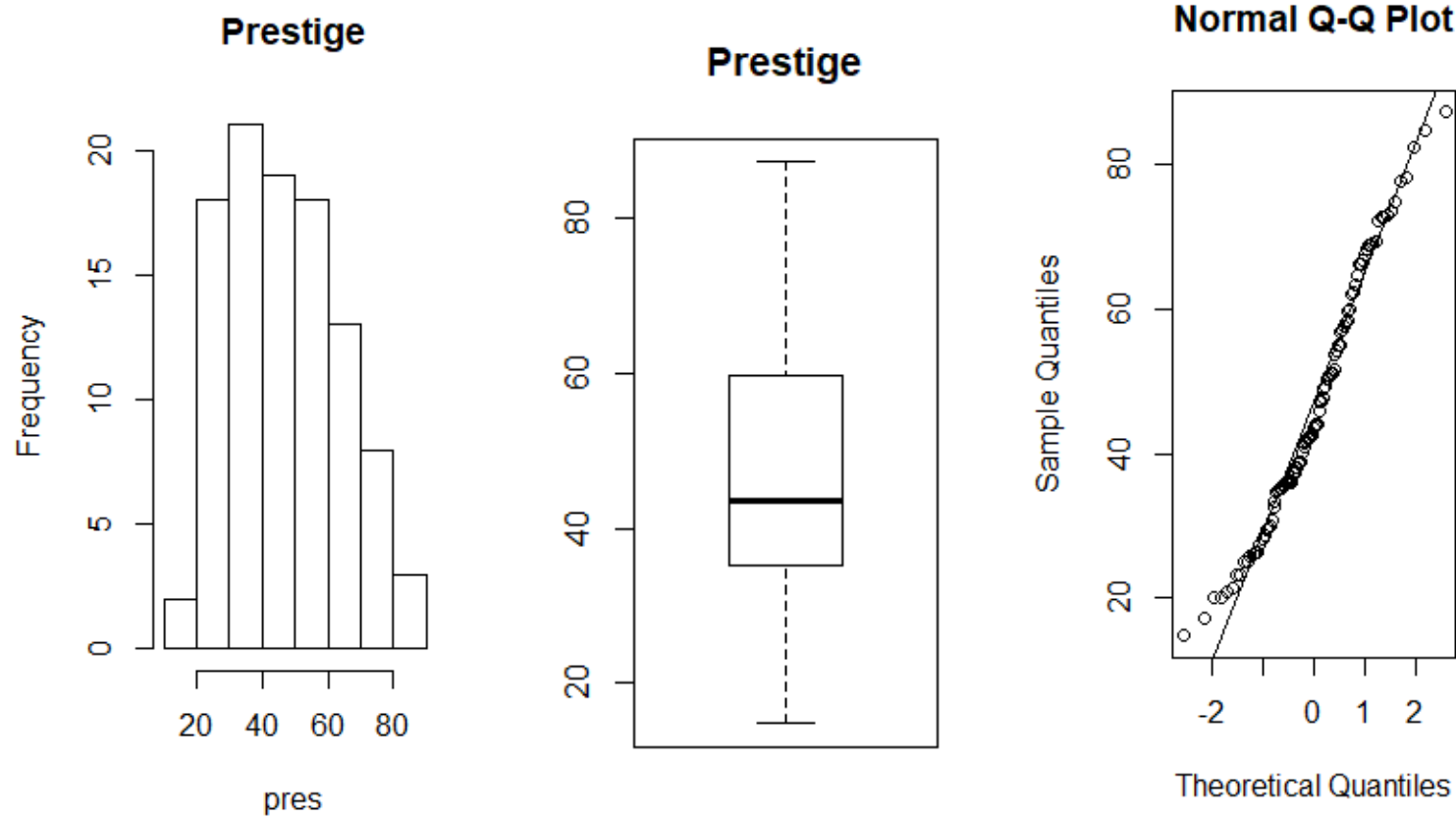
Για το εισόδημα:



Για τις γυναίκες επαγγελματίες:



Για την κοινωνικοοικονομική θέση:



Προχωρούμε στο *Shapiro – Wilk Test* με μηδενική υπόθεση ότι ακολουθείται η κανονική κατανομή και έχουμε τα εξής:

Για εκπαίδευση:  $W = 0.94958$ ,  $p - value = 0.0006773$

Για εισόδημα:  $W = 0.81505$ ,  $p - value = 5.634e - 10$

Για τις γυναίκες επαγγελματίες:  $W = 0.81579$ ,  $p - value = 5.957e - 10$

Για την κοινωνικοοικονομική θέση:  $W = 0.97198$ ,  $p - value = 0.02875$ .

Αυτό σημαίνει ότι δεν απορρίπτουμε τη μηδενική υπόθεση μόνο για τη μεταβλητή *prestige* σε επίπεδο στατιστικής σημαντικότητας 1%.

Άρα βρίσκουμε για τη μεταβλητή αυτή της μέση τιμή (46.83333), τη διασπορά (295.9943) και την τυπική απόκλιση (17.20449).

Για τις μεταβλητές που δεν ακολουθούν κανονική κατανομή μπορούμε να βρούμε το μέτρο θέσης *median* και τα ποσοστιαία *quantiles*. Παίρνουμε λοιπόν:

Για εκπαίδευση: $median = 10.54$ , $quantiles =$	0%	25%	50%	75%	100%
	6.3800	8.4450	10.5400	12.6475	15.9700

Για εισόδημα: $median = 5930.5$ , $quantiles =$	0%	25%	50%	75%	100%
	611.00	4106.00	5930.50	8187.25	25879.00

Για τις γυναίκες επαγγελματίες: $median = 13.6$ , $quantiles =$	0%	25%	50%	75%	100%
	0	3.5925	13.6	52.2025	97.5100

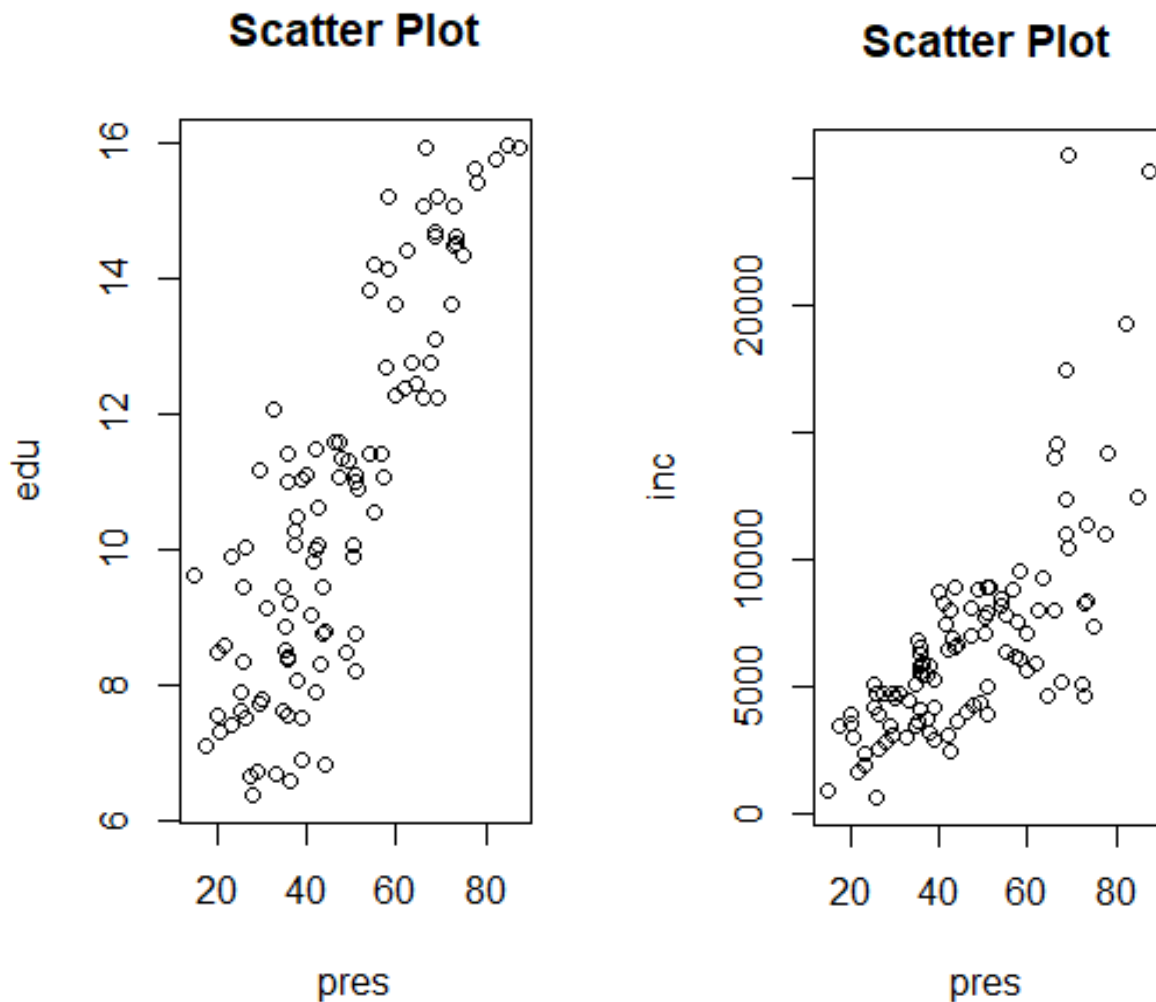
Παρατηρούμε λοιπόν ότι η μέση μόρφωση είναι 10.54 έτη, με την ελάχιστη στα 6.38 έτη και τη μέγιστη στα 15.97 έτη, το μέσο εισόδημα είναι 5930.5 δολάρια, με το ελάχιστο στα 611 και το μέγιστο στα 25879 δολάρια, ενώ το μέσο ποσοστό των εν ενεργεία γυναικών επαγγελματιών είναι 13.6%, με το ελάχιστο στο 0% και το μέγιστο στο 97.51%.

Το επάγγελμα με το μικρότερο εισόδημα είναι οι *babysitters*, ενώ με το υψηλότερο είναι οι *general managers*. Το επάγγελμα με το λιγότερο *prestige* είναι τα *newsboys*, ενώ με το μεγαλύτερο είναι οι *physicians*.



2ο Ερώτημα:

α) Παρουσιάζουμε τα δύο σημειογράμματα μεταξύ των μεταβλητών *prestige – education* και *prestige – income*:



Παρατηρούμε πως είναι γραμμική η σχέση εκπαίδευσης με κοινωνικοοικονομικής θέσης και όσο αυξάνεται η μία μεταβλητή, αυξάνεται και η άλλη. Από το σημειόγραμμα μεταξύ εισοδήματος και κοινωνικοοικονομικής θέσης, βλέπουμε πως ενώ μέχρι ένα σημείο υπάρχει έντονη γραμμική σχέση, όταν αυξηθεί λίγο ακόμα το *prestige*, το εισόδημα εκτοξεύεται.

β, γ) Χρησιμοποιώντας τη συνάρτηση *lm* της *R*, καταλήγουμε στα εξής:

Θεωρώντας ως εξαρτημένη μεταβλητή το *prestige* και ως ανεξάρτητη την *education* έχουμε κατάλοιπα:

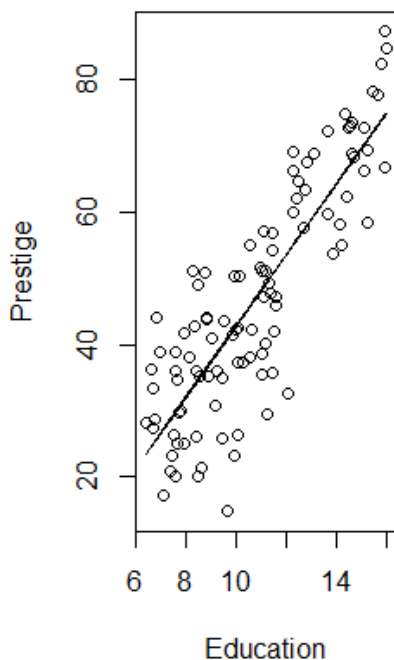
$Min$        $1Q$        $Median$        $3Q$        $Max$   
 $-26.0397$     $-6.5228$     $0.6611$     $6.7430$     $18.1636$ , οπότε συμπεραίνουμε πως υπάρχει αρκετά σημαντική συμμε-  
 τρικότητα στα ποσοστημόρια, αλλά όχι στα άκρα των καταλοίπων. Στην συνέχεια, έχουμε βρει τις τιμές των  
 $\beta_0$  και  $\beta_1$  μέσω εκτιμήσεων με το  $\hat{\beta}_0 = -10.732$  και το  $\hat{\beta}_1 = 5.361$  με αντίστοιχη τυπική απόκλιση 3.677 και 0.332.  
 Κάνουμε έλεγχο στατιστικής σημαντικότητας και προκύπτει ότι και τα δύο είναι σημαντικά για το δείγμα μας, το  
 πρώτο μέχρι 1 τοις χιλίοις και το δεύτερο μέχρι το 0. Οπότε καταλαβαίνουμε ότι η εκπαίδευση μπορεί καλύτερα  
 να προβλέψει την κοινωνικοοικονομική θέση από το ανάποδο. Παραθέτουμε τις τιμές των καταλοίπων και των  $\hat{Y}$ :

1	2	3	4	5
9.25087492	14.10762099	5.67357335	6.31075828	5.85594955
6	7	8	9	10
4.48785426	2.43633701	6.06002981	5.99203732	1.04873199
11	12	13	14	15
6.31070689	4.79318588	-9.60895705	-4.47909246	8.64977776
16	17	18	19	20
-10.34609058	8.49094016	-7.02443792	-12.56057709	5.79925487
21	22	23	24	25
9.71876461	-2.68317272	-4.01005421	12.37237339	-8.02040906
26	27	28	29	30
0.27347055	8.63544545	-5.02831259	9.81682728	-1.50696831
31	32	33	34	35
9.66635579	8.47984794	0.19522601	3.50354073	-5.40059093
36	37	38	39	40
-8.96450316	-0.55315394	-4.00775709	-2.46758905	-7.52726682
41	42	43	44	45
-21.38102979	-9.75210818	-2.59531071	-6.05205678	-7.51084298
46	47	48	49	50
-19.90984861	2.16541283	-14.84285049	-12.63767307	-0.51905490
51	52	53	54	55
-8.73458717	-16.64483922	-26.03966180	-19.20153390	-4.15419971
56	57	58	59	60
-1.62015206	2.70150060	3.46446986	3.73758837	-1.06121167
61	62	63	64	65
-14.63547874	8.96750436	-14.08192136	-7.76325180	-10.08385870
66	67	68	69	70
-9.80347123	18.16357829	-13.87156651	-1.57261228	9.21096388
71	72	73	74	75
-5.02512389	4.57487611	-5.84573079	8.16770995	3.39966606
76	77	78	79	80
-0.85927433	7.70264916	1.60060903	10.07383034	1.43978270
81	82	83	84	85
7.36347549	14.57069305	-7.23144988	4.72958205	5.40887235
86	87	88	89	90
6.88711689	2.27492751	3.11603851	7.69846610	17.65834947
91	92	93	94	95
12.53470807	11.55018895	-1.23647311	8.97587047	-3.08181857
96	97	98	99	100
11.05401221	14.11813004	5.99652877	-6.67977844	-8.03856464
101	102			
-0.67679534	0.09647737			

1	2	3	4	5	6	7	8
59.54913	54.99238	57.72643	50.48924	67.64405	73.11215	70.16366	72.03997
9	10	11	12	13	14	15	16
67.10796	67.75127	55.68929	55.20681	63.40896	66.67909	66.25022	65.44609
17	18	19	20	21	22	23	24
73.80906	65.12444	70.86058	67.00075	74.88124	62.28317	70.11005	74.82763
25	26	27	28	29	30	31	32
74.72041	68.12653	56.06455	39.92831	62.28317	70.80697	57.83364	48.72015
33	34	35	36	37	38	39	40
57.40477	50.59646	51.40059	50.86450	49.95315	46.30776	50.16759	38.42727
41	42	43	44	45	46	47	48
54.08103	48.45211	38.69531	43.25206	45.61084	49.30985	48.93459	50.54285
49	50	51	52	53	54	55	56
48.23767	42.01905	48.93459	43.14484	40.83966	42.50153	51.45420	48.72015
57	58	59	60	61	62	63	64
48.39850	40.03553	47.86241	30.76121	34.83548	45.93250	39.98192	28.56325
65	66	67	68	69	70	71	72
27.38386	29.90347	25.93642	35.37157	36.87261	29.68904	30.22512	30.22512
73	74	75	76	77	78	79	80
29.04573	25.13229	25.40033	43.35927	36.49735	34.29939	31.72617	34.46022
81	82	83	84	85	86	87	88
36.33652	36.22931	44.43145	23.47042	32.69113	43.41288	25.02507	37.78396
89	90	91	92	93	94	95	96
42.50153	33.44165	26.36529	24.64981	31.13647	33.92413	29.58182	55.04599
97	98	99	100	101	102		
34.78187	29.90347	31.77978	34.13856	42.87680	35.10352		

Χρησιμοποιώντας τα αποτελέσματα της συνάρτησης  $lm$ , κάνουμε πάλι σημειόγραμμα και παρατηρούμε ότι η γραμμή είναι ίδια με πριν, συμμετρικά, αφού άλλαξαν θέση οι μεταβλητές:

**Scatterplot F.R.L.**



Όμοια ακριβώς λειτουργούμε για το απλό γραμμικό μοντέλο των *prestige – income*.  
 όδημα εκτοξεύεται.

Χρησιμοποιώντας τη συνάρτηση *lm* της *R*, καταλήγουμε στα εξής:

Θεωρώντας ως εξαρτημένη μεταβλητή το *prestige* και ως ανεξάρτητη το *income* έχουμε κατάλοιπα:

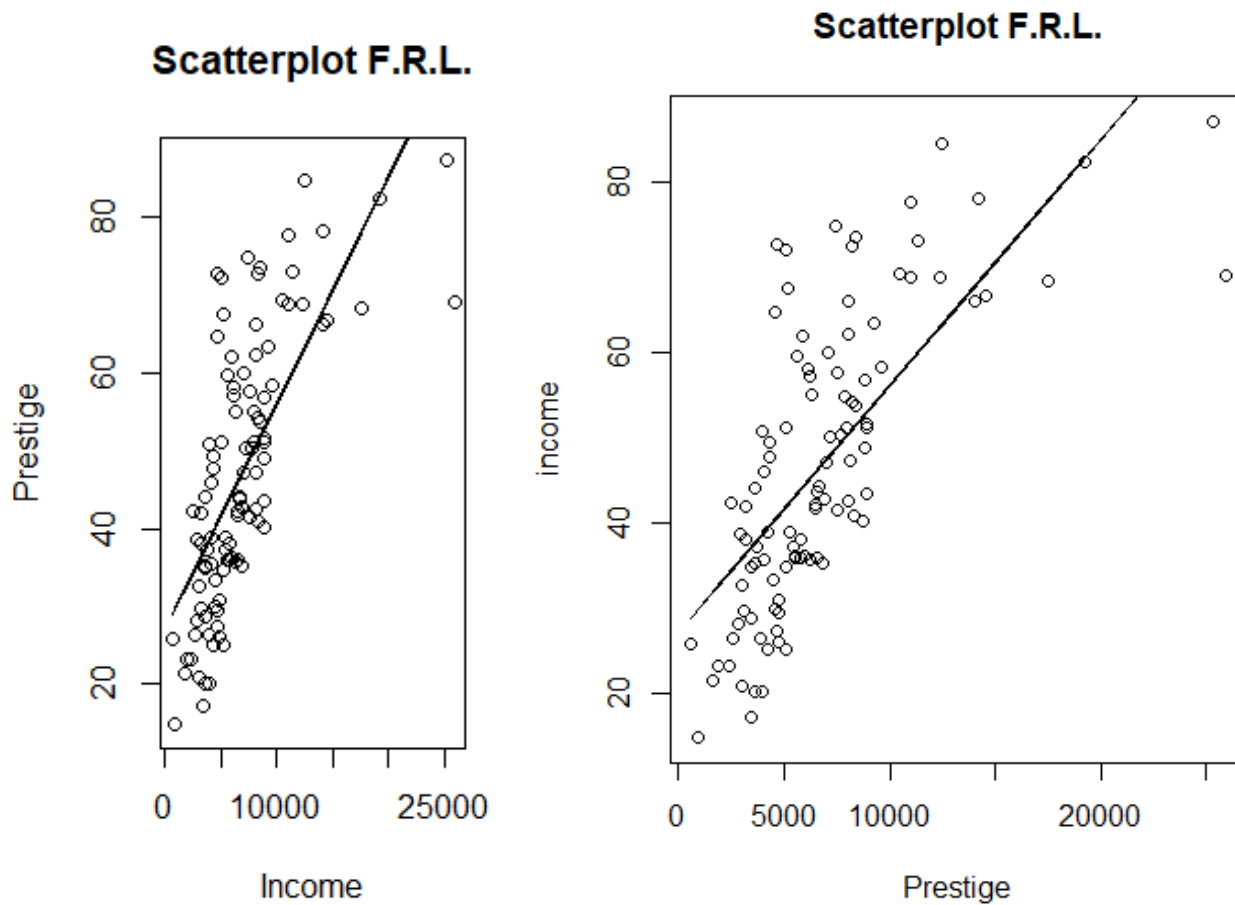
Min	1Q	Median	3Q	Max
-33.007	-8.378	-2.378	8.432	32.084

οπότε συμπεραίνουμε πως αφού η διάμεσος απέχει πολύ από το 0 δεν υπάρχει η συμμετρικότητα που είχαμε πριν, αλλά τα ποσοστημόρια και τα άκρα εμφανίζουν μεταξύ τους σημαντική συμμετρικότητα. Στην συνέχεια, έχουμε βρει τις τιμές των  $\beta_0$  και  $\beta_1$  μέσω εκτιμήσεων με το  $\hat{\beta}_0=27.141176368$  και το  $\hat{\beta}_1=0.002896799$  με αντίστοιχη τυπική απόκλιση  $2.268e+00$  και  $2.833e-04$ . Κάνουμε έλεγχο στατιστικής σημαντικότητας και προκύπτει ότι και τα δύο είναι σημαντικά για το δείγμα μας, με κοινή τιμή  $2e-16$ , δηλαδή κοντά στο 0. Οπότε καταλαβαίνουμε ότι η μία μεταβλητή προβλέπει την άλλη στο ίδιο επίπεδο. Παραθέτουμε τις τιμές των καταλοίπων και των  $\hat{Y}$ :

1	2	3	4	5	6
5.8804567	-33.0074429	9.4025982	3.9786987	22.0170199	18.5071285
7	8	9	10	11	12
21.5370558	9.9314566	13.0019391	9.7274060	17.7619148	12.4103181
13	14	15	16	17	18
2.2532904	11.7424869	26.3080256	9.6047039	-0.6422194	13.2535869
19	20	21	22	23	24
3.3698289	32.0844226	21.3067696	16.0977017	15.6859389	-13.2533705
25	26	27	28	29	30
-2.6127791	-9.4293688	24.1929921	-2.3365216	30.2083221	11.9394144
31	32	33	34	35	36
25.3534038	12.1073590	8.5532281	3.1876894	7.1673421	5.6396997
37	38	39	40	41	42
9.6635407	8.0674592	8.0156831	-10.0328374	-3.1779228	3.1552092
43	44	45	46	47	48
-7.0054368	-0.7723086	1.8020414	-11.4749014	9.3241941	-9.5722426
49	50	51	52	53	54
-3.3456331	-7.3150280	-12.3750733	-8.1554735	-15.0004380	-10.7065905
55	56	57	58	59	60
-3.3950507	-0.2955964	0.9118892	-9.4082053	-1.2966181	-6.4676027
61	62	63	64	65	66
-18.3255972	4.9639107	-3.0111207	-15.0315740	-19.8988632	-17.4175111
67	68	69	70	71	72
6.4057841	-10.4382758	-11.7132189	-0.4048362	-16.8133435	-7.2133435
73	74	75	76	77	78
-9.4161269	-6.7116552	-8.4365216	-7.9401323	-2.3091758	-10.2586631
79	80	81	82	83	84
-4.1037448	-8.0744765	-2.4818375	12.2396412	-5.7258352	-7.1883637
85	86	87	88	89	90
-5.8281277	0.8071210	-13.4445454	-10.3309585	2.3553997	-1.7647533
91	92	93	94	95	96
-3.5913153	-8.2032028	-10.4187159	-4.3102012	-11.9676612	-1.6890627
97	98	99	100	101	102
-3.8633653	-7.3531735	-14.2772562	-14.8096630	-3.6602928	-2.4188991

1	2	3	4	5	6	7
62.91954	102.10744	53.99740	52.82130	51.48298	59.09287	51.06294
8	9	10	11	12	13	14
68.16854	60.09806	59.07259	44.23809	47.58968	51.54671	50.45751
15	16	17	18	19	20	21
48.59197	45.49530	82.94222	44.84641	54.93017	40.71558	63.29323
22	23	24	25	26	27	28
43.50230	50.41406	100.45337	69.31278	77.82937	40.50701	37.23652
29	30	31	32	33	34	35
41.89168	57.36059	42.14660	45.09264	49.04677	50.91231	38.83266
36	37	38	39	40	41	42
36.26030	39.73646	34.23254	39.68432	40.93284	35.87792	35.54479
43	44	45	46	47	48	49
43.10544	37.97231	36.29796	40.87490	41.77581	45.27224	38.94563
50	51	52	53	54	55	56
48.81503	52.57507	34.65547	29.80044	34.00659	50.69505	47.39560
57	58	59	60	61	62	63
50.18811	52.90821	52.89662	36.16760	38.52560	49.93609	28.91112
64	65	66	67	68	69	70
35.83157	37.19886	37.51751	37.69422	31.93828	47.01322	39.30484
71	72	73	74	75	76	77
42.01334	42.01334	32.61613	40.01166	37.23652	50.44013	46.50918
78	79	80	81	82	83	84
46.15866	45.90374	43.97448	46.18184	38.56036	42.92584	35.38836
85	86	87	88	89	90	91
43.92813	49.49288	40.74455	51.23096	47.84460	52.86475	42.49132
92	93	94	95	96	97	98
44.40320	40.31872	47.21020	38.46766	67.78906	52.76337	43.25317
99	100	101	102			
39.37726	40.90966	45.86029	37.61890			

Χρησιμοποιώντας τα αποτελέσματα της συνάρτησης *lm*, κάνουμε πάλι σημειόγραμμα και παρατηρούμε ότι η γραμμή είναι ίδια με πριν, σχετικά συμμετρικά, αφού άλλαξαν θέση οι μεταβλητές:



δ) Για να κάνουμε έλεγχο υποθέσεων χρειάστηκε να εκτιμήσουμε τα  $\beta_0, \beta_1$  και το  $\Upsilon$ . Με αυτόν τον τρόπο υπολογίζουμε τα κατάλοιπα, τη διακύμανση και στη συνέχεια τη διακύμανση των  $\beta_0, \beta_1$ . Έχοντας ως υποθέσεις τις  $H_0 : \beta_1 = 0$  vs  $H_1 : \beta_1 > 0$  βρίσκουμε την παρατηρούμενη τιμή της ελεγχοσυνάρτησης  $t^* = 16.1478$ . Για επίπεδο στατιστικής σημαντικότητας  $\alpha=0.05$  με 100 βαθμούς ελευθερίας έχουμε ότι το ποσοστημόριο είναι 1.660234, πολύ μικρότερο από την παρατηρούμενη τιμή, άρα απορρίπτουμε την  $H_0$  και δεχόμαστε την εναλλακτική υπόθεση.

ε) Μέσω της συνάρτησης *lm* βρίσκουμε τα διαστήματα εμπιστοσύνης σε επίπεδο στατιστικής σημαντικότητας 10% για τους συντελεστές του μοντέλου *prestige – education* να κυμαίνονται από -16.83% έως -4.63% για το  $\beta_0$  και από 4.8% έως 5.91% για το  $\beta_1$ . Μέσω της συνάρτησης *lm* βρίσκουμε ομοίως τα διαστήματα εμπιστοσύνης σε

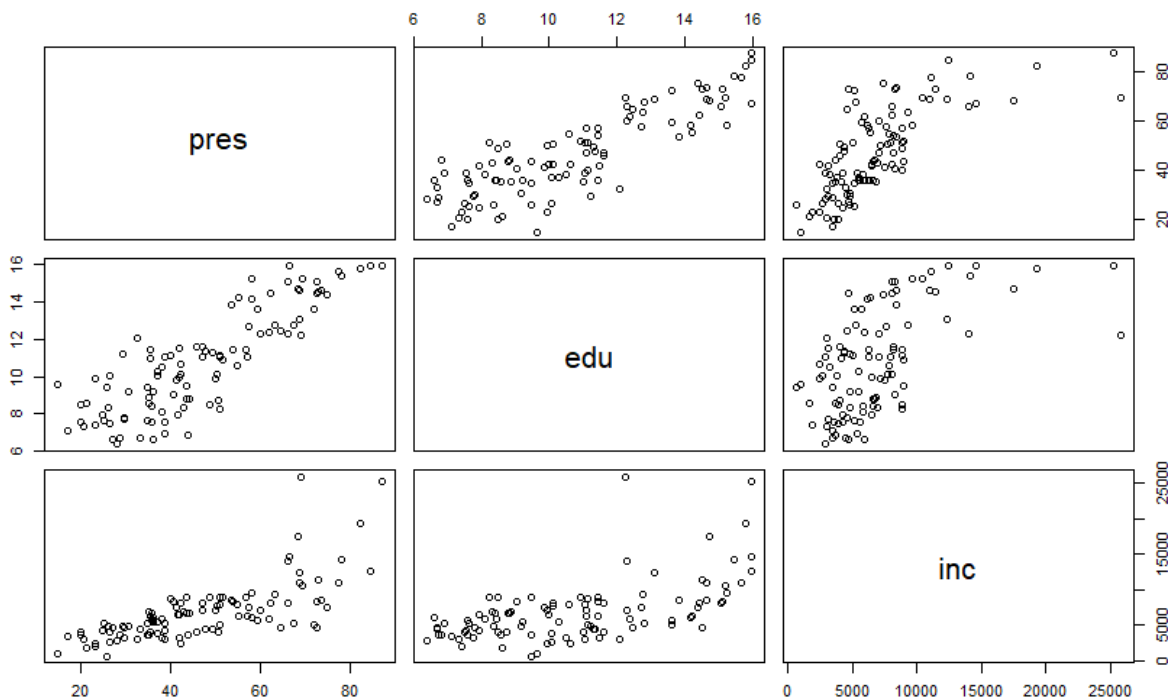
επίπεδο στατιστικής σημαντικότητας 10% για τους συντελεστές του μοντέλου *prestige – income* να κυμαίνονται από 23.37% έως 30.9% για το  $\beta_0$  και από 0.0024% έως 0.003367% για το  $\beta_1$ .

στ) Υπολογίζουμε μέσω του *summary* τον συντελεστή προσδιορισμού και τον προσαρμοσμένο συντελεστή προσδιορισμού για τα δύο μοντέλα κι έχουμε ότι για το *prestige – education*  $R^2 = 0.7228$ ,  $R^2_{adj} = 0.72$  και για το *prestige – income*  $R^2 = 0.5111$ ,  $R^2_{adj} = 0.5062$ . Στη γραμμική παλινδρόμηση του *prestige – education* το  $R^2$  είναι αρκετά μεγάλο ώστε να έχει πολλή πληροφορία, ενώ στο *prestige – income* είναι 0.5 περίπου ο συντελεστής προσδιορισμού. Αυτό σημαίνει ότι η μισή πληροφορία βρίσκεται σε σφάλματα οπότε χρειάζονται περισσότερα στοιχεία για καλό αποτέλεσμα, δηλαδή καλή προσαρμογή του μοντέλου.

ζ) Με βάση τα παραπάνω βρίσκουμε σημειακή εκτίμηση για την τιμή 42.87 σε διάστημα 41.32602 με 44.42757. Αντίστοιχα, το ατομικό διάστημα έχει ίδια πρόβλεψη, αλλά σε διάστημα 27.68385 με 58.06974.

3ο Ερώτημα:

α) Θέλουμε να οπτικοποιήσουμε τις σχέσεις των μεταβλητών που επιλέξαμε, οπότε κάνουμε το εξής σημειόγραμμα:



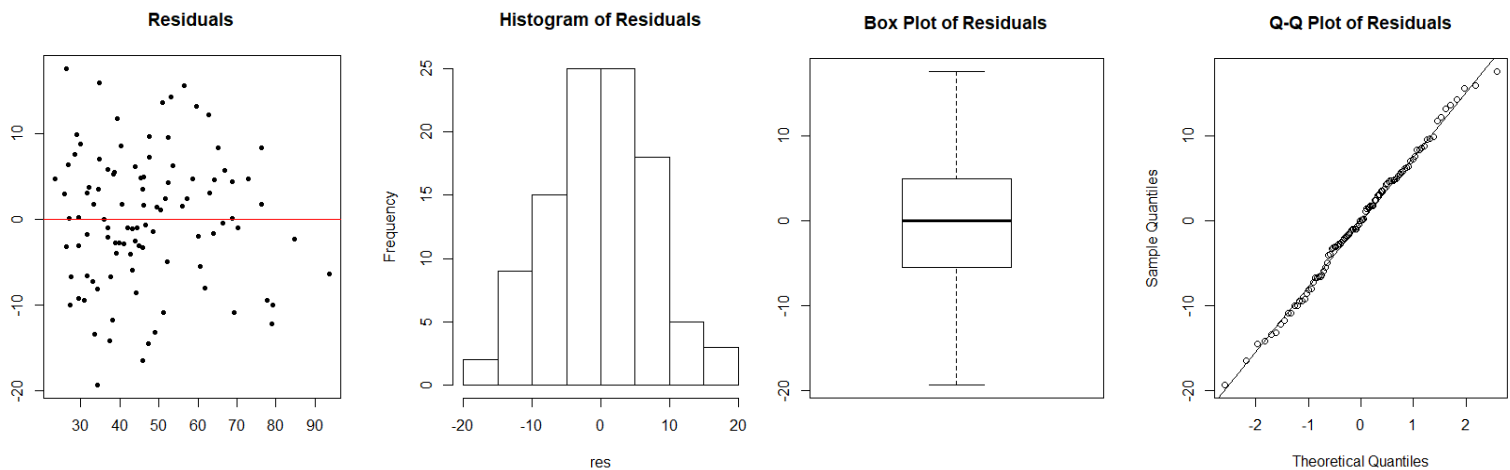
Δεδομένου ότι πήραμε ως εξαρτημένη μεταβλητή την κοινωνικοοικονομική θέση, εξετάζουμε την πρώτη γραμμή των σημειογραμμάτων. Παρατηρούμε θετική γραμμική σχέση πιο έντονη με την ανεξάρτητη μεταβλητή της εκπαίδευσης. Τα διαστήματα εμπιστοσύνης που εξάγουμε είναι:

	2.5 %	97.5 %
(Intercept)	-1.323493e+01	-0.460629799
edu	3.445127e+00	4.829761535
inc	9.162805e-04	0.001806051

Όσον αφορά τις παραμέτρους της παλινδρόμησης, παρατηρούμε από τα  $p$ -values, τα οποία βγήκαν  $< 2e-16$  και  $2.36e-08$  για την εκπαίδευση και το εισόδημα αντίστοιχα, ότι είναι στατιστικά σημαντικές και οι δύο μεταβλητές για το μοντέλο μας.

β) Μέσω της συνάρτησης *lm* υπολογίζουμε τον συντελεστή προσδιορισμού: 0.798 και τον προσαρμοσμένο: 0.7939. Αυτό σημαίνει ότι εξηγείται ένα 80% της μεταβλητότητας της κοινωνικοοικονομικής θέσης.

γ) Οι αρχικές υποθέσεις ελέγχονται μέσω του  $F$ -test. Κάνουμε τον έλεγχο για ολόκληρο το μοντέλο και καταλήγουμε σε τέτοιο  $p$ -value που δείχνει ότι το μοντέλο μας έχει νόημα όπως είναι. Κάνουμε τώρα γραφικό έλεγχο των καταλοίπων.



Όπως παρατηρούμε από τα παραπάνω, το ιστόγραμμα λόγω της ομοιομορφίας του στο κέντρο μας παραπέμπει σε κανονική κατανομή, αλλά για επιβεβαίωση προχωρούμε στο *boxplot* της παλινδρόμησης, όπου το κουτί του διαγράμματος σχεδόν ισαπέχει από τις άκρες και χωρίς ακραίες παρατηρήσεις. Σχεδιάζουμε στο τέλος το *q-qplot* που δείχνει κανονική κατανομή με ελάχιστες παρατηρήσεις να μην ακουμπούν τη γραμμή της κανονικής. Ολοκληρώνουμε τον έλεγχο με ένα *Shapiro test*, που δίνει  $p$ -value=0.9371. Αν θεωρήσουμε  $H_0$  την υπόθεση ότι το δείγμα ακολουθεί κανονική κατανομή και  $H_1$  την εναλλακτική υπόθεση ότι το δείγμα δεν ακολουθεί κανονική κατανομή, έχουμε ότι αν απορρίψουμε την  $H_0$  θα έχουμε  $\alpha$  γφάλμα της τάξης του 93%, οπότε δεν την απορρίπτουμε.



δ) Συγκρίνοντας τις τιμές των  $R^2$  και  $R^2_{adj}$  από τα μοντέλα έχουμε ότι στις απλές παλινδρομήσεις εξηγείται μικρότερο ποσό πληροφορίας σε σχέση με το πολλαπλό μοντέλο. Η διαφορά δεν είναι μεγάλη αφού για *prestige* και εκπαίδευση  $R^2=0.7228$  και για *prestige* και εισόδημα  $R^2=0.5111$ , ενώ στην πολλαπλή παλινδρόμηση  $R^2=0.798$ . Επομένως, θα προτιμήσουμε την ανάλυση με πολλαπλή παλινδρόμηση.

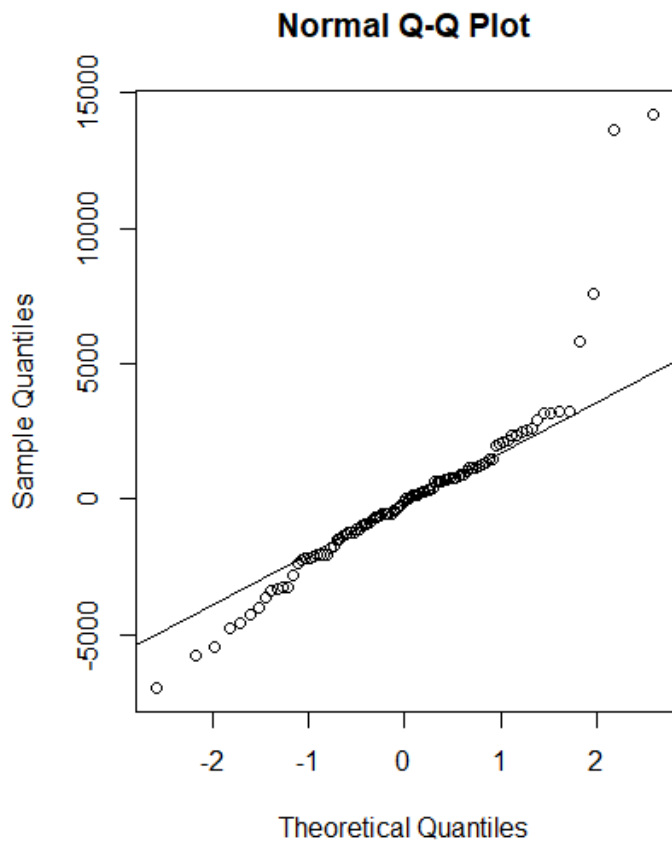
ε) Έχουμε βρει στο Ερώτημα 2 ότι σημειακή εκτίμηση 42.87 σε διάστημα 41.32602 με 44.42757. Αντίστοιχα, το ατομικό διάστημα έχει ίδια πρόβλεψη, αλλά σε διάστημα 27.68385 με 58.06974. Τώρα, βρίσκουμε σημειακή εκτίμηση 39.97133 σε διάστημα 38.42149 με 41.52116. Αντίστοιχα, το ατομικό διάστημα έχει ίδια πρόβλεψη, αλλά σε διάστημα 26.91113 με 53.03152. Βλέπουμε πάλι ότι η ατομική πρόβλεψη είναι πιο ευρεία για να εξασφαλίσουν επιτυχή πρόβλεψη.

στ) Για να κάνουμε έλεγχο υποθέσεων χρειάστηκε να εκτιμήσουμε τα  $\beta$  και το  $\Upsilon$ . Με αυτόν τον τρόπο υπολογίζουμε τα κατάλοιπα, τη διακύμανση και στη συνέχεια τη διακύμανση των  $\beta$ . Έχοντας ως υποθέσεις τις  $H_0 : \beta_{edu} = 4$  vs  $H_1 : \beta_{edu} > 4$  βρίσκουμε την παρατηρούμενη τιμή της ελεγχοσυνάρτησης  $t^* = -3.280814e + 12$ . Για επίπεδο στατιστικής σημαντικότητας  $\alpha=0.05$  με 100 βαθμούς ελευθερίας έχουμε ότι το ποσοστημόριο είναι 1.660715, πολύ μικρότερο από την παρατηρούμενη τιμή, άρα δεν μπορούμε να απορρίψουμε την  $H_0$ .

#### 4ο Ερώτημα:

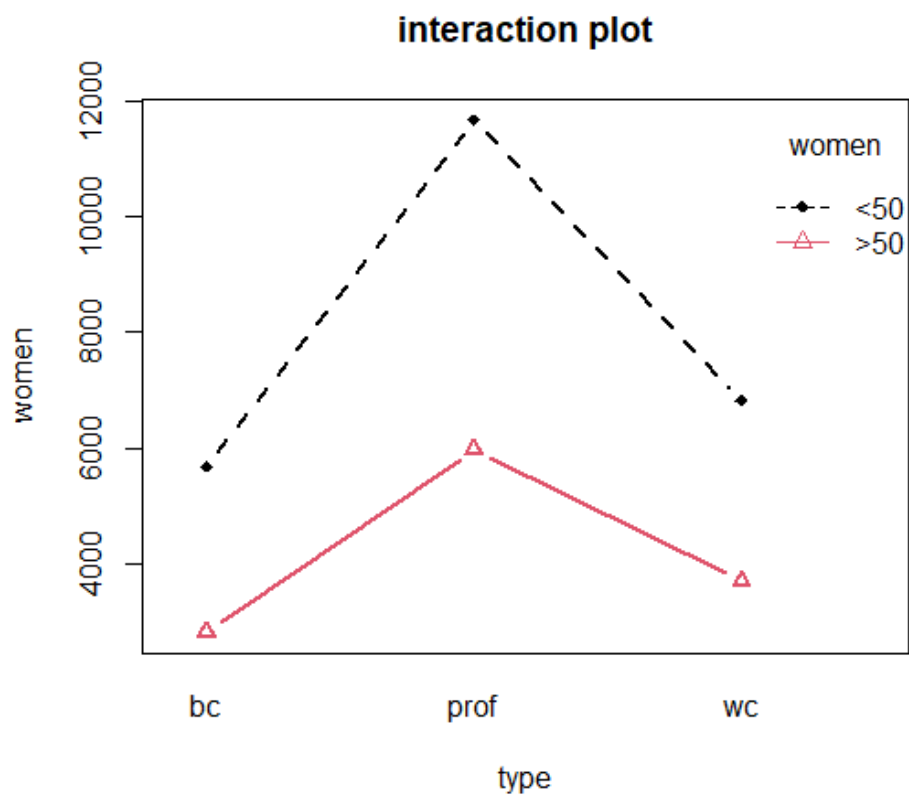
α, β) Υποθέσεις του μοντέλου ANOVA: ι) η κατανομή των παρατηρήσεων σε κάθε επίπεδο είναι κανονική για το επίπεδο  $i$ , ιι) η κανονική κατανομή σε κάθε επίπεδο έχει την ίδια διασπορά  $\sigma^2$ . Ισοδύναμα, οι τυχαίοι όροι  $\epsilon_{i,j}$  είναι ταυτοτικά κατανομημένοι για κάθε  $i, j$ , ιιι) οι παρατηρήσεις σε κάθε επίπεδο του παράγοντα είναι ανεξάρτητες και ταυτοτικά κατανομημένες και είναι ανεξάρτητες από τις παρατηρήσεις στα άλλα επίπεδα.

Κάνουμε τις κατάλληλες αλλαγές ώστε η μεταβλητή *women* να παίρνει τις τιμές που ζητήθηκαν και διαχωρίζουμε εξαρτημένη και ανεξάρτητες μεταβλητές. Φτάνουμε έτσι σε ένα μοντέλο ANOVA 2 παραγόντων, όπου και οι δύο ανεξάρτητες μεταβλητές είναι στατιστικά σημαντικές, αλλά η αλληλεπίδρασή τους όχι, με  $p - value=0.271$ . Βγάζοντας την αλληλεπίδραση από το μοντέλο, καταλήγουμε πως μπορεί διαφορετικά δρουν οι γυναίκες επαγγελματίες από τις χειρώναχτες μεταξύ τους, αλλά με τον ίδιο τρόπο σε σχέση με το εισόδημα. Έτσι, καταλήγουμε στην τελική μορφή του μοντέλου, που περιέχει τις μεταβλητές *women, type*, αλλά όχι την αλληλεπίδρασή τους. γ) Θέλουμε να ελέγξουμε τα κατάλοιπα οπότε κάνουμε το  $q - plot$ :



Παρατηρούμε ότι στις ουρές χαλάει η κανονικότητα οπότε ελέγχουμε με *Shapiro test* και βρίσκουμε  $p\text{-value} = 5.38e - 09$ , τόσο μικρό που απορρίπτουμε την κανονικότητα και στηρίζομαστε στο μεγάλο δείγμα για την αξιοπιστία του μοντέλου. Προηγουμένως, είχαμε οδηγηθεί σε κατάλοιπα που ακολουθούσαν κανονική κατανομή σε αντίθεση με τώρα.

δ) Κάνουμε το *interaction plot* και βλέπουμε ότι οι δύο μεταβλητές έχουν παρόμοια μορφή, άρα έχουν την ίδια επίδραση περίπου στο μοντέλο. Μάλιστα, αυτός είναι ένας πολύ καλός τρόπος να καταλάβουμε ότι σωστά βγάλαμε τον όρο της αλληλεπίδρασης προηγουμένως από το μοντέλο, καθώς παρατηρούμε μικρή διαφορά:



## ΚΩΔΙΚΑΣ R:

```
1 data<-read.csv(file=prestige, header=TRUE)
2 head(prestige)
3 profs<-prestige$Professions
4 edu<-prestige$education
5 inc<-prestige$income
6 women<-prestige$women
7 pres<-prestige$prestige
8 type<-prestige$type
9 type_freq=table(type)
10 type_freq
11 profs_freq=table(profs)
12 barplot(type_freq,main="Type of worker", xlab="Number of workers", col="black")
13 pie(type_freq, main="Type of worker", col=c("grey10", "grey20", "grey30"))
14 par(mfrow=c(1,2))
15 hist(edu,main="Education")
16 hist(inc, main="Income")
17 hist(women, main="women")
18 hist(pres, main="Prestige")
19 boxplot(edu,main="Education")
20 boxplot(inc, main="Income")
21 boxplot(women, main="women")
22 boxplot(pres, main="Prestige")
23 qqnorm(edu)
24 qqline(edu)
25 qqnorm(inc)
26 qqline(inc)
27 qqnorm(women)
28 qqline(women)
29 qqnorm(pres)
30 qqline(pres)
31 shapiro.test(edu)
32 shapiro.test(inc)
33 shapiro.test(women)
34 shapiro.test(pres)
35 mu=mean(pres)
36 mu
37 variance=var(pres)
38 variance
39 sd=sd(pres)
40 sd
41 median(edu)
42 quantile(edu)
43 median(inc)
44 quantile(inc)
45 median(women)
46 quantile(women)
47 min(inc)
48 incom<-which.min(inc)
49 profs[incom]
50 max(inc)
51 inco<-which.max(inc)
52 profs[inco]
53 min(pres)
54 prestig<-which.min(pres)
55 profs[prestig]
56 max(pres)
```

```

57 prest<-which.max(pres)
58 profs[prest]
59 plot(pres,edu,
60      main="Scatter Plot",
61      preslab="Prestige",
62      edulab="Education")
63 plot(pres,inc,
64      main="Scatter Plot",
65      preslab="Prestige",
66      inclab="Income")
67 y<-prestige$pres
68 x<-prestige$edu
69 reg<-lm(y~x)
70 summary(reg)
71 reg$coefficients
72 reg$residuals
73 reg$fitted.values
74 plot(x,y,
75      main="Scatterplot F.R.L.",
76      xlab="Education",
77      ylab="Prestige")
78 lines(x,reg$fitted.values)
79 x<-prestige$inc
80 reg<-lm(y~x)
81 summary(reg)
82 reg$coefficients
83 reg$residuals
84 reg$fitted.values
85 plot(x,y,
86      main="Scatterplot F.R.L.",
87      xlab="Income",
88      ylab="Prestige")
89 lines(x,reg$fitted.values)
90 reg1=lm(pres~inc)
91 reg1$coefficients
92 reg1$residuals
93 reg1$fitted.values
94 z=pres
95 w=inc
96 plot(w,z,
97      main="Scatterplot F.R.L.",
98      xlab="Prestige",
99      ylab="income")
100 lines(w,reg1$fitted.values)
101 shapiro.test(edu)
102 shapiro.test(inc)
103 shapiro.test(pres)
104 x<-edu
105 t<-sum((x-mean(x))^2)
106 b1<-sum((x-mean(x))*(y-mean(y)))/t
107 b0<-mean(y)-b1*mean(x)
108 yhat<-b0+b1*x
109 ehat<-y-yhat
110 sigma2hat<-sum(ehat^2)/(100)
111 sigmahat<-sqrt(sigma2hat)

```

```

112 sigma2_b0<-sigma2hat*(1/100+mean(x)^2/t)
113 std_b0=sqrt(sigma2_b0)
114 sigma2_b1<-sigma2hat/t
115 std_b1=sqrt(sigma2_b1)
116 c = 0
117 t_stat<-(b1-c)/std_b1
118 t_stat
119 a<-0.05
120 qt_a<-qt(1-a, 100)
121 qt_a
122 t<-sum((w-mean(w))^2)
123 b1<-sum((w-mean(w))*(y-mean(y)))/t
124 b0<-mean(y)-b1*mean(w)
125 yhat<-b0+b1*w
126 ehat<-y-yhat
127 n<-102
128 sigma2hat<-sum(ehat^2)/(n-2)
129 sigma2hat<-sqrt(sigma2hat)
130 sigma2_b0<-sigma2hat*(1/n+mean(w)^2/t)
131 std_b0=sqrt(sigma2_b0)
132 sigma2_b1<-sigma2hat/t
133 std_b1=sqrt(sigma2_b1)
134 c = 0
135 t_stat<-(b1-c)/std_b1
136 t_stat
137 a<-0.05
138 qt_a<-qt(1-a, n-2)
139 qt_a
140 reg=lm(pres~edu)
141 reg1=lm(pres~inc)
142 confint(reg, level=0.90)
143 confint(reg1, level=0.90)
144 summary(reg)
145 summary(reg1)
146 regg=lm( pres ~ edu + inc )
147 summary(regg)
148 confint(regg)
149 futurex<-10
150 futurew<-4000
151 future<-data.frame(edu=futurex, inc=futurew)
152 predict(reg, newdata=future, level=0.90, interval="confidence")
153 predict(reg, newdata=future, level=0.90, interval="prediction")
154 head(prestige)
155 pres<-prestige$prestige
156 edu<-prestige$education
157 inc<-prestige$income
158 pairs(cbind(pres,edu,inc))
159 reg <- lm(pres ~ edu + inc)
160 summary(reg)
161 confint(reg)
162 coef=reg$coefficients
163 res =reg$residuals
164 fitted=reg$fitted.values
165 plot(fitted,res,main = "Residuals",xlab = "",ylab = "",
166      pch=20)

```

```

167 lines(x=c(0,100),y=c(0,0),col="red"))
168 hist(res,main="Histogram of Residuals")
169 boxplot(res,main="Box Plot of Residuals")
170 n=length(profs)
171 ones=rep(1,n)
172 x=cbind(ones,edu,inc)
173 betas=solve(t(X)%*%X)%*%t(X)%*%pres
174 ehat=y-X%*%betas
175 qqnorm(res,main="Q-Q Plot of Residuals")
176 qqline(ehat)
177 shapiro.test(res)
178 simple.edu<-lm(pres~edu)
179 simple.inc<-lm(pres~inc)
180 summary(simple.edu)
181 summary(simple.inc)
182 summary(reg)
183 res.edu =simple.edu$residuals
184 res.inc =simple.inc$residuals
185 res =reg$residuals
186 edu.avg<-10
187 inc.avg<-4000
188 newjob<-data.frame(edu=edu.avg, inc= inc.avg)
189 predict(reg, newdata=newjob, level=0.90, interval="confidence")
190 predict(reg, newdata=newjob, level=0.90, interval="prediction")
191 betas<-solve(t(X)%*%X)%*%t(X)%*%y
192 betas
193 yhat<-X%*%betas
194 ehat<-y-X%*%betas
195 ehat
196 sigma2.hat<-t(ehat)%*%ehat/(n-2)
197 sigma.hat<-sqrt(sigma2.hat)
198 var.betas<-as.vector(sigma2.hat)*solve(t(X)%*%X)
199 std.betas<-sqrt(diag(var.betas))
200 var.betas
201 std.betas
202 betas
203 c<-4
204 t_stat<-(betas[2]-c)/std.betas[2]
205 t_stat
206 a<- 0.05
207 qt_a<-qt(1-a, n-5)
208 qt_a
209 install.packages("gplots")
210 library(gplots)
211 head(prestige)
212 women<-prestige$women
213 women[women<50]<-0
214 women[women >=50]<-1
215 y<-prestige$income
216 n<-length(y)
217 y
218 n
219 type<-prestige$type
220 women
221 type
222 women<-factor(women)
223 levels(women)<- c("<50", ">50")
224 two.way.aov<-aov(y~ women*type)
225 summary(two.way.aov)
226 two.way.aov<-aov(y~ women+type)
227 summary(two.way.aov)
228 resid.twoway<-two.way.aov$residuals
229 qqnorm(resid.twoway)
230 qqline(resid.twoway)
231 shapiro.test(resid.twoway)
232 interaction.plot(type,women,y,type="b",col=c(1:3),leg.bty="n",lwd=2,
233                 pch=c(18,24),xlab="type",ylab="women",main="interaction plot")

```